

LLM alignment fine-tuning: DPO vs IPO

Tri Nguyen

nguyetr9@oregonstate.edu

March 27, 2025

1 “Notations”

- a LLM **generates** responses: perform certain decoding techniques given a language model. Popular choices: beam search, top-k sampling, top-p sample.
-

2 What is LLM fine-tuning and why

Large language models (LLMs) are becoming significantly popular and bring broad impacts across various aspects of our life. One of the key factors attributed to this success is the capability of a model to perform on very wide range of tasks: you can just input a task description and the relevant information and feed it to the model as a prompt. This ability is commonly referred to LLM being able to acquire various skills, or being able to follow instruction, or can be seen as LLM’s generalization power **careful here**.

This amazing performance can be largely attributed to the fine-tuning steps, where the pretrained model is twisted to learn to produce outputs that are more aligned with human (or the creator’s goals), whether it is to follow instruction, to produce “better response”, including being helpful, harmless, factual, or refuse to response to malicious requests. The very first work that shakes the whole NLP/AI community is from OpenAI, where they kind of lay out the roadmap to do fine-tuning:

1. Supervised fine-tuning. Assume access to demonstration dataset: $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ where \mathbf{y}_n is the ideal response given the prompt \mathbf{x} . **Example here pls**. The objective in the step is to twist the model to produce more response like in the datasets. This is supervised learning objective and can be accomplished as any next-token predicting task.

The challenge in this step is not about the method but about the dataset. It is considerably hard to collect the ideal response that is varying over the whole prompt space. Even OpenAI can only collect $[[\mathbf{x}]]$ pair of samples to do this. Therefore the need of second step, which is the focus of this article.

2. Preference learning. There are some variants, but many works considers a using pairwise preference dataset: $\{\mathbf{x}_n, \mathbf{y}_n^{(1)}, \mathbf{y}_n^{(2)}, c_n\}$, where given a prompt \mathbf{x}_n , and 2 possible responses $\mathbf{y}_n^{(1)}, \mathbf{y}_n^{(2)}$, the preference label $c_n \in \{1, 2\}$ indicates which response is preferred over the others.

Example: This type of dataset is arguably easier to collect: one can use pretrained LLMs to generate 2 different responses given the prompt \mathbf{x} , then ask a human annotator to choose which one is better. The particular details such as where to collect \mathbf{x} , how to define preference are depending on specific tasks. And although many works tend to refer the problem as human preference fine-tuning, preference can be anything, not necessarily based on human.

Using this dataset, the popular work Reinforcement Learning with Human Feedback (RLHF) proposed to formulate the problem of fine-tuning LLM as:

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\mathbb{E}_{\mathbf{y} \sim \pi} [s(\mathbf{x}, \mathbf{y}) - \beta D_{\text{kl}}(\pi \parallel \pi_{\text{ref}})] \right], \quad (1)$$

where the *score function* $s(\mathbf{x}, \mathbf{y})$ outputs a scalar value indicating how strongly aligned the response \mathbf{y} is given prompt \mathbf{x} . The objective in English is: learn a policy (a language model) π that behaves not so much differently from π_{ref} and π generates responses \mathbf{y} with highest possible $s(\mathbf{x}, \mathbf{y})$. π_{ref} is a reference policy which can be seen as a good initialization to start with. In practice, π_{ref} is the trained policy obtained from step 1.

Having the overall view of the fine-tuning pipeline, we are ready to dive what the existing works proposed.

3 RLHF and DPO

3.1 RLHF

While the title didn't mention RLHF, we start from RLHF to see the cleverness in the development of DPO. The very first obstacle in optimizing (1) is that score function s is **unknown**. RLHF proposed to solve this as an independent problem. If we assume the preference label follow the Bradley-Terry (BT) model, we have

$$\Pr(c = 1 \mid \mathbf{x}, \mathbf{y}^{(1)}, \mathbf{y}^{(2)}) = \frac{\exp(s(\mathbf{x}, \mathbf{y}^{(1)}))}{\exp(s(\mathbf{x}, \mathbf{y}^{(1)})) + \exp(s(\mathbf{x}, \mathbf{y}^{(2)}))} = \sigma(s(\mathbf{x}, \mathbf{y}^{(1)}) - s(\mathbf{x}, \mathbf{y}^{(2)})) \quad (2)$$

As the score function gives a higher score for better aligned response, it is pretty reasonable to employ the BT model.

Pairwise Generative Model. Under this model, the problem of learning s can be cast as a somewhat special binary classification problem: Given a sample $(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)$, we wish to learn $h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)$ to predict the true binary label c .

$$\Pr(c = 1 \mid \mathbf{y}^a, \mathbf{y}^b, \mathbf{x}) = \sigma(s^{\sharp}(\mathbf{y}^a, \mathbf{x}) - s^{\sharp}(\mathbf{y}^b, \mathbf{x})) \quad (3)$$

There are some specific details that a general classification problem does not apply here:

1. Firstly, because of the underlying model given in (2), h_{θ} need to be parameterized as

$$h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2) \triangleq \sigma(s_{\theta}(\mathbf{x}, \mathbf{y}^1) - s_{\theta}(\mathbf{x}, \mathbf{y}^2))$$

where s_{θ} can be parameterized using a LLM. The argument for the use of LLM is that LLM is so large, it has the potential to work on any task, including giving a score to a pair of prompt, response. Then maximum likelihood estimation gives us

$$\underset{\theta}{\text{minimize}} \quad - \mathbb{E}_{\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2} [\mathbb{I}[c = 1] \log \sigma(h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)) + \mathbb{I}[c = 2] \log (1 - \sigma(h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)))] \quad (4)$$

2. Secondly,

Remark 3.1. If we obtain a sample $(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)$ with label 1, it is equivalent to obtain the sample $(\mathbf{x}, \mathbf{y}^2, \mathbf{y}^1)$ with label 2.

Proof.

$$\begin{aligned} \Pr(c = 1 \mid \mathbf{y}^1, \mathbf{y}^2, \mathbf{x}) &= 1 - \Pr(c = 2 \mid \mathbf{y}^1, \mathbf{y}^2, \mathbf{x}) \\ &= 1 - \sigma(s^{\sharp}(\mathbf{x}, \mathbf{y}^2) - s^{\sharp}(\mathbf{x}, \mathbf{y}^1)) \\ &= 1 - \Pr(c = 1 \mid \mathbf{y}^2, \mathbf{y}^1, \mathbf{x}) \\ &= \Pr(c = 2 \mid \mathbf{y}^2, \mathbf{y}^1, \mathbf{x}) \end{aligned}$$

■

Well, this information is already encoded in the parameterization of $h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2)$, i.e., $h_{\theta}(\mathbf{x}, \mathbf{y}^1, \mathbf{y}^2) = 1 - h_{\theta}(\mathbf{x}, \mathbf{y}^2, \mathbf{y}^1)$. We will come back to this point later since not every method inherently embed this structure in their parameterization.

Because of this, we can rewrite the optimization as:

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}^w, \mathbf{y}^l} \left[-\log \sigma(h_{\theta}(\mathbf{x}, \mathbf{y}^w, \mathbf{y}^l)) \right], \quad (5)$$

After solving (5), the alignment fine-tuning is performed by plugging the trained s_{θ^*} into (1) and employ any off-the-shelf RL techniques, e.g, PPO [1]. So that is the story of RLHF.

The two-step approach is inherently quite complex and computation intensive. That's where DPO comes to play.

3.2 DPO

An optimal solution of (1) can be derived as follows:

$$\begin{aligned} & \arg \max_{\pi} \quad \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [s(\mathbf{y}, \mathbf{x}) - \beta D_{\text{kl}}(\pi \parallel \pi_{\text{ref}})] \right] \\ &= \arg \min_{\pi} \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} \left[\log \frac{\pi(\mathbf{y} | \mathbf{x})}{\exp(\beta^{-1} s(\mathbf{y}, \mathbf{x})) \pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right] \right] \\ &= \arg \min_{\pi} \mathbb{E}_{\mathbf{x}} \left[D_{\text{kl}}(\pi \parallel \frac{1}{Z(\mathbf{x})} \exp(\beta^{-1} s(\mathbf{y}, \mathbf{x})) \pi_{\text{ref}}(\mathbf{y} | \mathbf{x})) \right] \\ &= \frac{1}{Z(\mathbf{x})} \exp(\beta^{-1} s(\mathbf{y}, \mathbf{x})) \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \end{aligned}$$

And hence,

$$\begin{aligned} \pi^*(\mathbf{y} | \mathbf{x}) &= \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp \left(\beta^{-1} s^{\dagger}(\mathbf{y}, \mathbf{x}) \right), \\ \Leftrightarrow s^{\dagger}(\mathbf{y}, \mathbf{x}) &= \beta \log \frac{\pi^*(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} + \beta \log Z(\mathbf{x}) \end{aligned}$$

where $Z(\mathbf{x})$ is an intractable normalizing factor.

Using the above identity, the generative model in (3) is equivalent to

$$\Pr(c = 1 | \mathbf{y}^1, \mathbf{y}^2, \mathbf{x}) = \sigma \left(\beta \log \frac{\pi^*(\mathbf{y}_1 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1 | \mathbf{x})} - \beta \log \frac{\pi^*(\mathbf{y}_2 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2 | \mathbf{x})} \right) \quad (6)$$

The new generative model is equivalent to (3), while offering a specific parameterization of the score function, i.e.,

$$s_{\theta}(\mathbf{y}, \mathbf{x}) = \beta \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})}.$$

Using this parameterization, the optimal solution θ^* of the sigmoid criterion (5) directly gives an optimal policy π_{θ^*} with respect to the fine-tuning objective (1). This solution eliminates the need of RL step as in RLHF.

is it truly equivalent?

3.3 A unified view

Pairwise Preference Generative model.

$$\Pr(c = 1 | \mathbf{y}^1, \mathbf{y}^2, \mathbf{x}) = \sigma \left(q^{\dagger}(\mathbf{y}_1, \mathbf{x}) - q^{\dagger}(\mathbf{y}_2, \mathbf{x}) \right),$$

The choice of $q^{\dagger}(\mathbf{y}_1, \mathbf{x})$ depends on the interest of learning underlying model.

- RLHF wishes to learn the score function: $q^\natural(\mathbf{y}_1, \mathbf{x}) = s^\natural(\mathbf{y}_1, \mathbf{x})$
- DPO wishes to learn the optimal policy: $q^\natural(\mathbf{y}, \mathbf{x}) = \beta \log \frac{\pi^\star(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})}$

In any case, after deciding the choice for q^\natural , we can use maximum likelihood to estimate q^\natural with q_θ .

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell} [-\log \sigma(q_\theta(\mathbf{y}_w, \mathbf{x}) - q_\theta(\mathbf{y}_\ell, \mathbf{x}))] \quad (7)$$

4 IPO

This paper is kind of hard to read as the notation is a bit messed up. I have mixed feeling about it: the idea is quite nice after I figured out all the technical details and fought against my intuition during all the way. However the presentation of their development is quite *cryptic*.

Let’s start with their empirical loss:

$$\underset{\theta}{\text{minimize}} \quad (q_\theta(\mathbf{y}_w, \mathbf{x}) - q_\theta(\mathbf{y}_\ell, \mathbf{x}) - 0.5)^2,$$

where q_θ is defined as in DPO.

This criterion is very counter-intuitive. Specifically,

- Turning a classification-based to a regression-based objective. This is ahhhhhh. But this is just because of my view, nothing wrong here.
- Why regressing toward 0.5? What special about 0.5? And why regressing all possible samples toward 0.5? Does it mean that we treat preference evenly among all samples? While in comparison to (7), the gap $q_\theta(\mathbf{y}_w, \mathbf{x}) - q_\theta(\mathbf{y}_\ell, \mathbf{x})$ to be as large as possible.

Well, I have some answered, but not all sadly. So intuition does not help here. Let’s dive into the technical development to see if that helps.

4.1 IPO’s Loss Derivation

In this section, I attempted to (i) follow and demystify IPO’s derivation, make the whole thing a bit more rigorous to identify the key step, and (ii) make a slightly more general formulation.

Note that IPO’s goal is to perform the fine-tuning task via (1) as RLHF’s and DPO’s. What distinguish them is their proposed score function:

$$s^\natural(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{y}' \sim \mu} [v(\mathbf{y}, \mathbf{y}', \mathbf{x})], \quad (8)$$

where the scalar-valued function $v(\mathbf{y}, \mathbf{y}', \mathbf{x}) \in \mathbb{R}$ denotes the degree of preference, i.e., $v(\mathbf{y}_1, \mathbf{y}'_1, \mathbf{x}) \succ v(\mathbf{y}_2, \mathbf{y}'_2, \mathbf{x})$ implies that the preference of \mathbf{y}_1 over \mathbf{y}'_1 is “stronger” compared to the preference comparison of \mathbf{y}_2 vs \mathbf{y}'_2 . Here, μ is just some arbitrary policy.

With that physical meaning, $v(\mathbf{y}, \mathbf{y}', \mathbf{x})$ is any function satisfying:

$$v(\mathbf{y}, \mathbf{y}', \mathbf{x}) \in [\alpha_1, \alpha_2], \quad \forall \mathbf{y}, \mathbf{y}', \mathbf{x} \quad (9a)$$

$$v(\mathbf{y}, \mathbf{y}, \mathbf{x}) = 0.5(\alpha_2 + \alpha_1), \quad \forall \mathbf{y}, \mathbf{x} \quad (9b)$$

$$v(\mathbf{y}, \mathbf{y}', \mathbf{x}) = \alpha_1 + \alpha_2 - v(\mathbf{y}', \mathbf{y}, \mathbf{x}) \quad (9c)$$

Condition (9a) is to avoid DPO’s issue in which over-fitting leads to policy deviating arbitrarily far away from π regardless of β . Condition (9b) and (9c) is to enforce the physical meaning of pairwise preference. With that score function, IPO’s generative model can be seen as follows.

IPO’s Pairwise Generative Model. The binary label $c \in \{1, 2\}$ for sample $\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}$ is

$$\Pr(c = 1) \propto v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \quad (10)$$

Here, binary RV $c = 1$ indicates $\mathbf{y}_1 \succ \mathbf{y}_2$. The likelihood of this even happening is proportional to $v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$, while the exact probability is unknown. From modeling perspective, the lack of specification in (10) relative to the BT model makes it more general.

Similar to DPO, the optimal policy to IPO's objective in (1) is

$$\pi^*(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} \mid \mathbf{x}) \exp \left(\beta^{-1} \mathbb{E}_{\mathbf{y}' \sim \mu} [v(\mathbf{y}, \mathbf{y}', \mathbf{x})] \right),$$

which enables the trick originally made by DPO by comparing 2 responses:

$$q^*(\mathbf{y}_1, \mathbf{x}) - q^*(\mathbf{y}_2, \mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})], \quad (11)$$

where $q^*(\mathbf{y}, \mathbf{x}) = \beta \log \pi^*(\mathbf{y} \mid \mathbf{x}) - \beta \log \pi_{\text{ref}}(\mathbf{y} \mid \mathbf{x})$. Let $q_\theta(\mathbf{y}, \mathbf{x}) = \beta \log \pi_\theta(\mathbf{y} \mid \mathbf{x}) - \beta \log \pi_{\text{ref}}(\mathbf{y} \mid \mathbf{x})$. All we want is to enforce the equality (11) for all $\mathbf{y}_1, \mathbf{y}_2$ with respect to q_θ , and one way to realize it is using squared loss

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[\left(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}) - \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right)^2 \right]$$

Expanding the squared term, the optimization problem reduces to

$$\min_{\theta} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))^2 - 2[q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})] \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right] \quad (12)$$

The issue lays in the unknown factor in the second term. Now comes to the IPO's peculiar derivations. The second term can be written as

$$\begin{aligned} & \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[[q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})] \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right] \\ &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})) (v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}))] \\ &= \underbrace{\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x}) v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})]}_{(*)} \\ &\quad - \underbrace{\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})]}_{(**)} \end{aligned}$$

For (*), notice that the two terms are essentially the same under expectation when $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}$ are i.i.d.

$$\begin{aligned} (*) &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x}) v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\ &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_2, \mathbf{y}} [q_\theta(\mathbf{y}_2, \mathbf{x}) v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\ &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_1, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] \\ &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \\ &= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x}) v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \end{aligned} \quad (13)$$

This trick is pretty cool ha ;), I also used this one in one of my proof :))) [x]. Note that the particular name of the variables under expectation might not matter as long as they share the same distribution. For example,

$$\mathbb{E}_{x \sim p} [f(x)] + \mathbb{E}_{y \sim p} [f(y)] = \mathbb{E}_{x \sim p} [f(x)] + \mathbb{E}_{x \sim p} [f(x)] = 2 \mathbb{E}_{x \sim p} [f(x)].$$

Similar trick can be applied for $(**)$ as well:

$$\begin{aligned}
(**) &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] \\
&= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= 2 \mathbb{E}_{\mathbf{y}_1 \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})] \mathbb{E}_{\mathbf{y}_2, \mathbf{y}} [v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= (\alpha_1 + \alpha_2) \mathbb{E}_{\mathbf{y}_1 \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})]
\end{aligned} \tag{14}$$

where the last equality holds because $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}$ are independent and $\mathbb{E}_{\mathbf{y}_2, \mathbf{y} \sim \mu} [v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] = 0.5(\alpha_1 + \alpha_2)$, which in turn can be derived from conditions (9b) and (9c). Combining (13), (14) gives

$$\begin{aligned}
(*) - (**) &= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] - (\alpha_1 + \alpha_2) \mathbb{E}_{\mathbf{y}_1} [q_\theta(\mathbf{y}_1, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [2q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - (v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) + v(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x}))q_\theta(\mathbf{y}_1, \mathbf{x})] \quad (\text{by (9c)}) \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x})q_\theta(\mathbf{y}_1, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})q_\theta(\mathbf{y}_2, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})]
\end{aligned}$$

The expression in (12) becomes

$$\begin{aligned}
&\arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))^2] - 2 \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \\
&= \arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}))^2] \\
&\approx \arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} (q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}) - \mathbb{I}[c=1]\alpha_2 - \mathbb{I}[c=2]\alpha_1)^2,
\end{aligned} \tag{15}$$

where $c \in \{1, 2\}$ is the pairwise preference label. The approximation is very crude: using $\{\alpha_1, \alpha_2\}$ to approximate $v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$. In a more general view, this is a non-parametric estimation of $v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$ under the generative model (10) using only 1 sample c . Note that this estimation is made *independently* to the alignment learning, i.e, optimizing (15). It is also the reason why IPO claims that they **even don't need to learn the score function**.

A further variance-reduced improvement made in IPO is to consider the symmetrical property of the sample: if we obtain the triple $(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$ with label c , then it is legitimate to assume another triple $(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x})$ with label $2/c$. With that in mind, we arrive at the final criterion:

$$\begin{aligned}
&\arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(h_\theta(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) - \mathbb{I}[c=1]\alpha_2 - \mathbb{I}[c=2]\alpha_1)^2 + (h_\theta(\mathbf{x}, \mathbf{y}_2, \mathbf{y}_1) - \mathbb{I}[c=1]\alpha_1 - \mathbb{I}[c=2]\alpha_2)^2] \\
&= \arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(h_\theta(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell) - \alpha_2)^2 + (h_\theta(\mathbf{x}, \mathbf{y}_{i,l}, \mathbf{y}_{i,w}) - \alpha_1)^2] \\
&= \arg \min_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(h_\theta(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell) - 0.5(\alpha_1 + \alpha_2))^2]
\end{aligned} \tag{16}$$

IPO's particular choice. In IPO, they settle with $v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) = \Pr(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x})$, and hence $\alpha_1 = 0, \alpha_2 = 1$. While it is looking sensible, any other arbitrary choice should be as valid. Anyhow, it leads to their criterion

$$\text{minimize}_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} [(h_\theta(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell) - 0.5)^2]$$

Effect of α_1, α_2 . Firstly, the final criterion (16) reveals that only the sum $\alpha_1 + \alpha_2$ matters, not the absolute values of α_1, α_2 . Large $\alpha_1 + \alpha_2$ means a broader preference model. However, it also leads to higher estimation error.

4.2 IPO’s loss discussion

- There is no score learning. The chosen score function has no parameter to estimate. This is contrast to RLHF/DPO’s approach where we known functional form of the score function governed by parameter θ .
- The seemingly disconnection between the score function and the preference model. However, in reality, they are connected via: μ and the oracle preference probability $\Pr(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x}) = f^\natural(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$, where f^\natural is any function satisfying $f^\natural(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \in [0, 1]$, $f^\natural(\mathbf{y}, \mathbf{y}, \mathbf{x}) = 0.5$, and $f^\natural(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) + f^\natural(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x}) = 1$. Okay, we are getting some structure here. Under this view, the preference data directly dictates the score function.
- There is no learning, there is no overfitting. We only use approximation
- In terms of preference model, anything is feasible, including contradicting preferences, i.e, $\mathbf{y}_1 \succ \mathbf{y}_2, \mathbf{y}_2 \succ \mathbf{y}_3, \mathbf{y}_3 \succ \mathbf{y}_1$. The BT model does not allow this. Therefore, if this happen because of noisy labels, IPO would just adapt to it while RLHF/DPO won’t.
- Although the preference model is very liberal, i.e, allow for contradicting preferences, the objective implies that there is a total ordering over responses \mathbf{y} .
- The score depends on μ , which is not necessarily a good feature.
- There could be this situation: $\Pr(\mathbf{y}_1 \succ \mathbf{y}_2)$ and $s(\mathbf{y}_1, \mathbf{x}) < s(\mathbf{y}_2, \mathbf{x})$.
- But the $f_\theta(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) = q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})$ has structure.
- Their motivation is to address DPO’s weakness. To me, it is more like a by-produce that they archive that goal. Let’s see.

5 Experiment Design

We want to show a single point: IPO’s approximation could be terribly wrong, and hence, it could not reach the optimal solution. In the meantime, DPO attain optimal solution easily. But can we design generative models such that the preference labels for both IPO and DPO are the same? The problem is the score functions are not the same.

Let assume bandit, 3 actions y_0, y_1, y_2 . For DPO, assume

$$s_{\text{DPO}}^\natural(\mathbf{y}) = [10.000, 9.593, 7.800] \quad (17)$$

$$s_{\text{IPO}}^\natural(y) = \mathbb{E}_{y'}[\Pr(y \succ y')] \quad (18)$$

There are 2 factors: score function, and preference generative model.

Method	Score function	Preference generative model
DPO	$s^\natural(\mathbf{y})$	$\Pr(c = 1) = \sigma(s^\natural(\mathbf{y}_1) - s^\natural(\mathbf{y}_2))$
IPO	$\mathbb{E}_{\mathbf{y}' \sim \mu}[v^\natural(\mathbf{y}, \mathbf{y}')]]$	$\Pr(c = 1) = v^\natural(\mathbf{y}_1, \mathbf{y}_2)$

Score function of IPO has structure, and it implies this interesting property: total score is a constant. Not sure how to feel about it!